

Protein databases

Henrik Nielsen

Background- Nucleotide databases

GenBank, <http://www.ncbi.nlm.nih.gov/Genbank/>

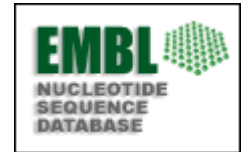
National Center for Biotechnology Information (NCBI), National Library of Medicine (NLM), National Institutes of Health (NIH), USA.



EMBL, <http://www.ebi.ac.uk/embl/>

European Bioinformatics Institute (EBI), England

(Established in 1980 by the European Molecular Biology Laboratory, Heidelberg, Tyskland)



DDBJ, <http://www.ddbj.nig.ac.jp/>

National Institute of Genetics, Japan



Together they form

International Nucleotide Sequence Database Collaboration,
<http://www.insdc.org/>



Protein databases

Swiss-Prot, <http://www.expasy.org/sprot/>

Established in 1986 in Switzerland

ExPASy (Expert Protein Analysis System)

Swiss Institute of Bioinformatics (SIB) and European Bioinformatics Institute (EBI)

PIR, <http://pir.georgetown.edu/>

Established in 1984

National Biomedical Research Foundation, Georgetown University, USA

In 2002 merged into:

UniProt, <http://www.uniprot.org/>

A collaboration between SIB, EBI and Georgetown University.



UniProt Knowledgebase (UniProtKB)

UniProt Reference Clusters (UniRef)

UniProt Archive (UniParc)

UniProt Knowledgebase Release 2011_02 (08-Feb-11)
consists of:

UniProtKB/Swiss-Prot: Annotated manually (*curated*)

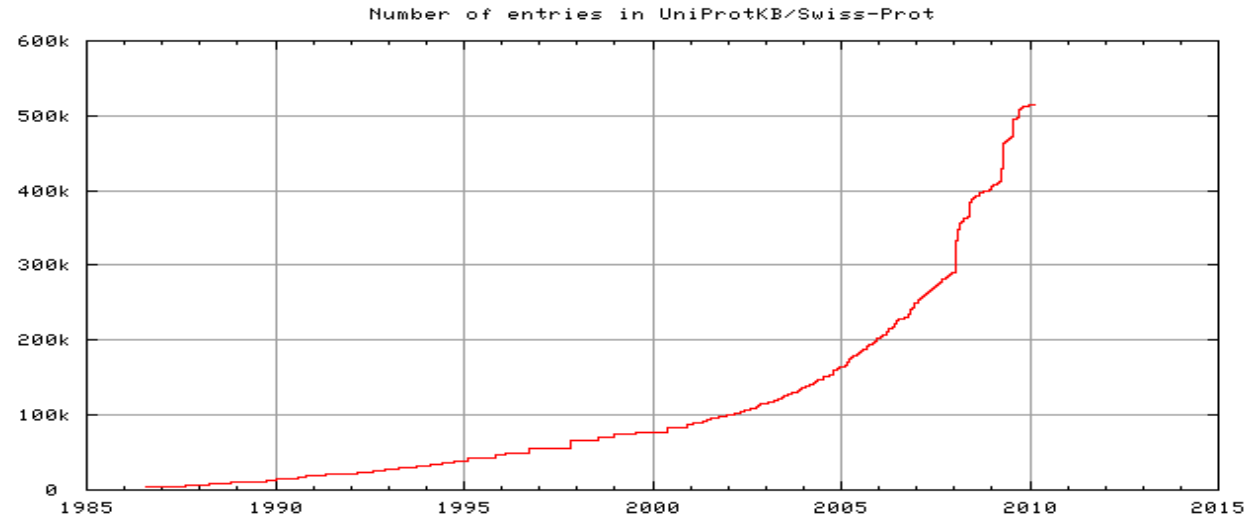
525,207 entries

UniProtKB/TrEMBL: Computer annotated

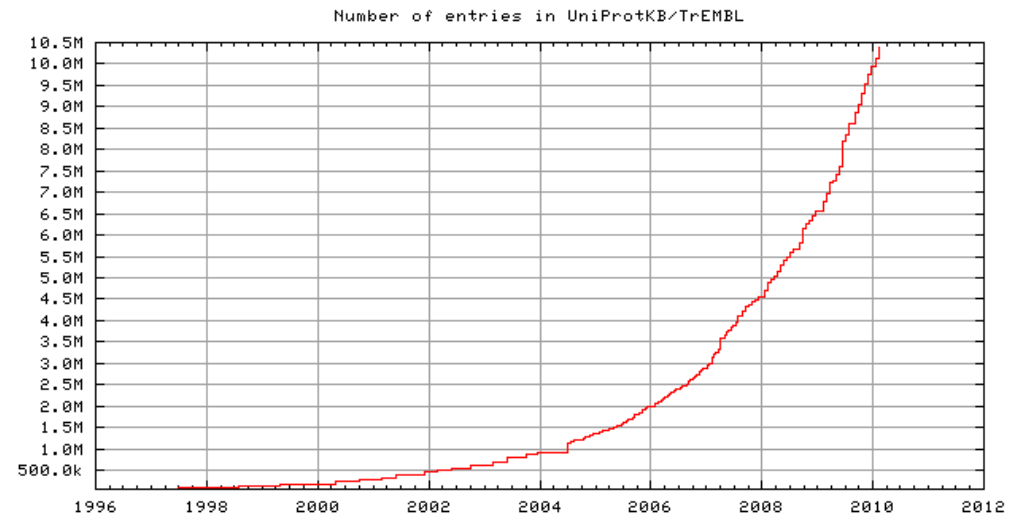
13,499,622 entries

Growth of UniProt

Swiss-Prot



TrEMBL



Content of UniProt Knowledgebase

- Amino acid sequences
 - Functional and structural annotations
 - Function / activity
 - Secondary structure
 - Subcellular location
 - Mutations, phenotypes
 - Post-translational modifications
 - Origin
 - organism: Species, subspecies; classification
 - tissue
 - References
 - Cross references
-

Amino acid sequences

From where do you get amino acid sequences?

- Translation of nucleotide sequences (GenBank/EMBL/DDBJ)
 - Amino acids sequencing: *Edman degradation*
 - Mass spectrometry
 - 3D-structures
-

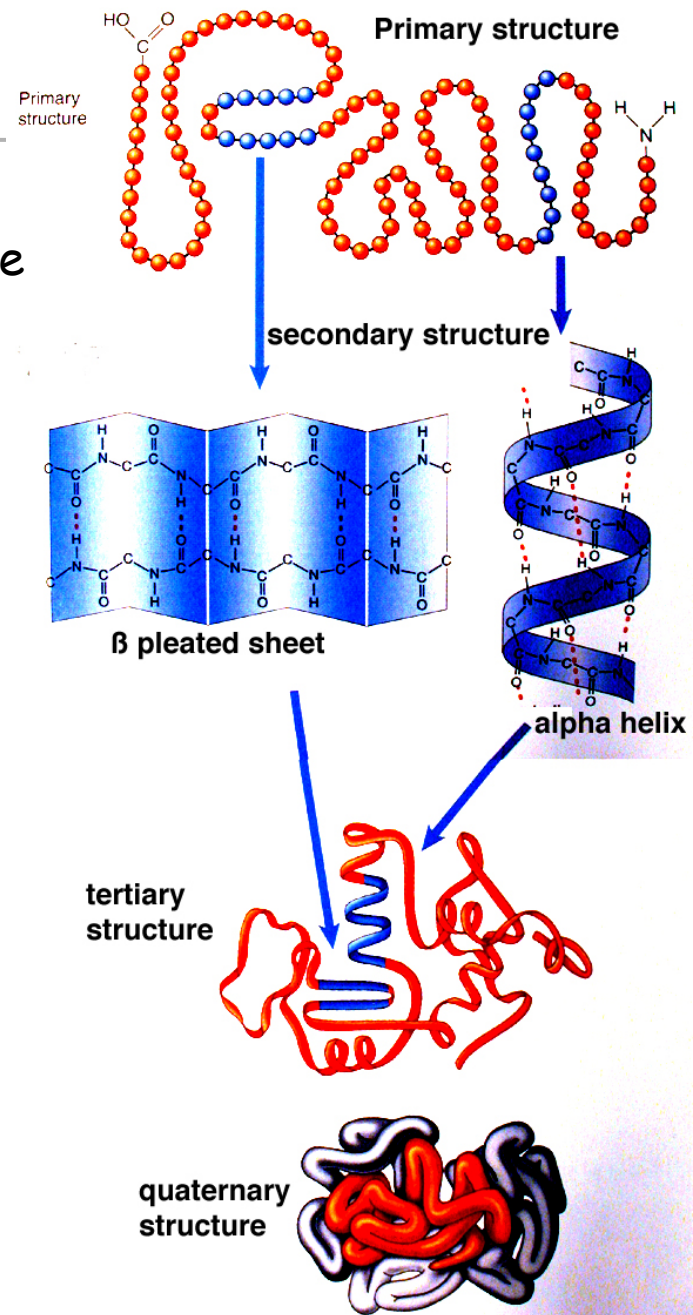
Protein structure

Primary structure: Amino acid sequence

Secondary structure:
"Backbone" hydrogen bonding
Alpha helix / Beta sheet / Turn

Tertiary structure: Fold, 3D coordinates

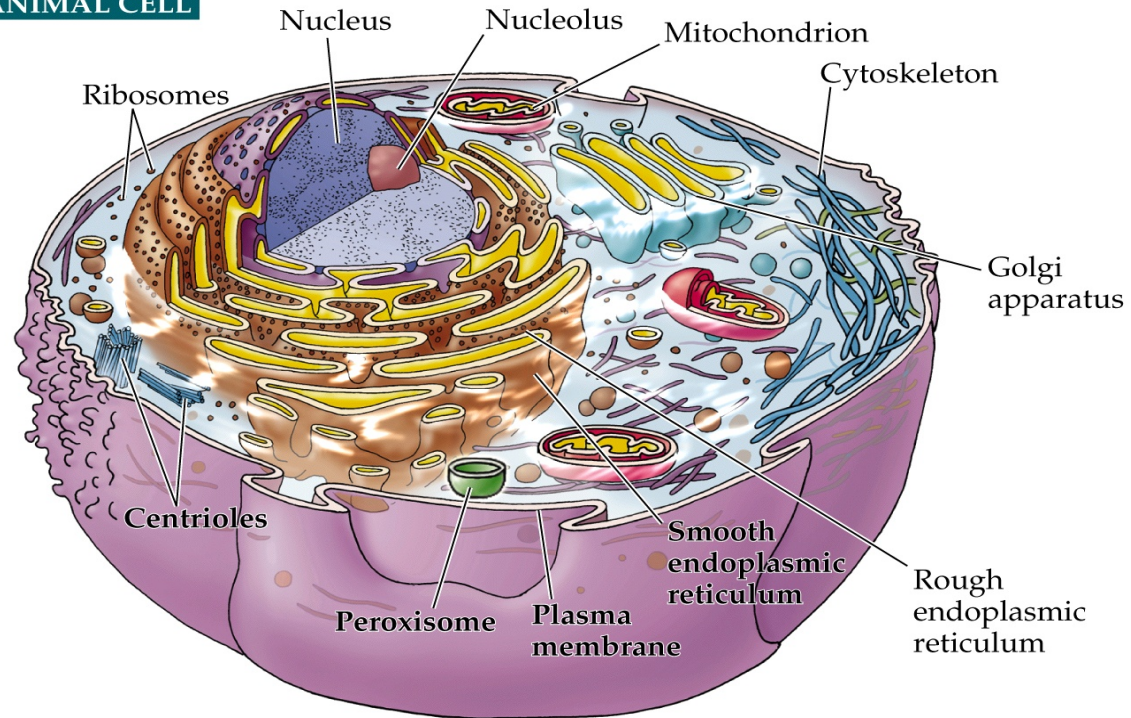
Quaternary structure: subunits



Subcellular location

An animal cell:

AN ANIMAL CELL



Post-translational modifications

- Cleavage of signal peptide, transit peptide or pro-peptide
 - Phosphorylation
 - Glycosylation
 - Lipid anchors
 - Disulfide bond
 - Prosthetic groups (*e.g.* metal ions)
-

Content of UniProt Knowledgebase

<http://www.uniprot.org/uniprot/Q9ULV8>

- Name, entry data etc
- Organism
- Functional annotations (comments)
- Sequence
- References
- Cross references
 - 3D structure - PDB
 - EMBL
 - ..

Evidence

- 3 types of *non-experimental qualifiers* in
Sequence annotation and General comment:
- *Potential*: Predicted using sequence analysis
 - *Probable*: Uncertain experimental evidence
 - *By similarity*: Predicted using sequence similarity
-

Cross references

Other databases (there are ~100 in total):

- Nucleotide sequences
 - 3D structure
 - Protein-protein interactions
 - Enzymatic activities and pathways
 - Gene expression (microarrays and 2D-PAGE)
 - Ontologies
 - Families and domains
 - Organism specific databases
-

The genetic code

		Second letter				
		U	C	A	G	
First letter	U	<div>UUU</div> <div>UUC</div> Phenylalanine <div>UUA</div> <div>UUG</div> Leucine	<div>UCU</div> <div>UCC</div> <div>UCA</div> <div>UCG</div> Serine	<div>UAU</div> <div>UAC</div> Tyrosine <div>UAA</div> <div>UAG</div> Stop codon Stop codon	<div>UGU</div> <div>UGC</div> Cysteine <div>UGA</div> <div>UGG</div> Stop codon Tryptophan	U C A G
	C	<div>CUU</div> <div>CUC</div> <div>CUA</div> <div>CUG</div> Leucine	<div>CCU</div> <div>CCC</div> <div>CCA</div> <div>CCG</div> Proline	<div>CAU</div> <div>CAC</div> Histidine <div>CAA</div> <div>CAG</div> Glutamine	<div>CGU</div> <div>CGC</div> <div>CGA</div> <div>CGG</div> Arginine	U C A G
	A	<div>AUU</div> <div>AUC</div> <div>AUA</div> Isoleucine <div>AUG</div> Methionine; start codon	<div>ACU</div> <div>ACC</div> <div>ACA</div> <div>ACG</div> Threonine	<div>AAU</div> <div>AAC</div> Asparagine <div>AAA</div> <div>AAG</div> Lysine	<div>AGU</div> <div>AGC</div> Serine <div>AGA</div> <div>AGG</div> Arginine	U C A G
	G	<div>GUU</div> <div>GUC</div> <div>GUA</div> <div>GUG</div> Valine	<div>GCU</div> <div>GCC</div> <div>GCA</div> <div>GCG</div> Alanine	<div>GAU</div> <div>GAC</div> Aspartic acid <div>GAA</div> <div>GAG</div> Glutamic acid	<div>GGU</div> <div>GGC</div> <div>GGA</div> <div>GGG</div> Glycine	U C A G

- Degenerate (*redundant*) but not ambiguous
- *Almost* universal (deviations found in mitochondria)

Reading Frames 1

A piece of an mRNA-strand:

5' aug ccc aag cug aa u agc gua gag ggg uuu uca uca uuu gag gac gau gua uaa 3'

can be divided into triplets (*codons*) in three ways:

1	aug	ccc	aag	cug	aa u	agc	gua	gag	ggg	uuu	uca	uca	uuu	gag	gac	gau	gua	uaa
	M	P	K	L	N	S	V	E	G	F	S	S	F	E	D	D	V	*
2	ugc	cca	agc	uga	a u a	g c g	uag	agg	gg u	uu u	ca u	ca u	u g	agg	ac g	au g	ua u	
	C	P	S	*	I	A	*	R	G	F	H	H	L	R	T	M	Y	
3	gcc	caa	gcu	gaa	uag	c g u	aga	ggg	gu u	uuc	auc	au u	uga	gga	cga	ugu	aua	
	A	Q	A	E	*	R	R	G	V	F	I	I	*	G	R	C	I	

Each possible set of triplets is called a *reading frame*.

Reading Frames 2

Since there are two strands in DNA, there are *six* possible reading frames in a piece of DNA (three in each direction):

3	A	Q	A	E	*	R	R	G	V	F	I	I	*	G	R	C	I		
2	C	P	S	*	I	A	*	R	G	F	H	H	L	R	T	<u>M</u>	<u>Y</u>		
1	<u>M</u>	<u>P</u>	<u>K</u>	<u>L</u>	<u>N</u>	<u>S</u>	<u>V</u>	<u>E</u>	<u>G</u>	<u>F</u>	<u>S</u>	<u>S</u>	<u>F</u>	<u>E</u>	<u>D</u>	<u>D</u>	<u>V</u>	*	
5'	ATGCCCAAGCTGAATAGCGTAGAGGGGTTTTTCATCATTTGAGGACGATGTATAA																	3'	
3'	TACGGGTTCGACTTATCGCATCTCCCCAAAAGTAGTAAACTCCTGCTACATATT																	5'	
	H	G	L	Q	I	A	Y	L	P	K	*	*	K	L	V	I	Y	L	-1
		G	L	S	F	L	T	S	P	N	E	D	N	S	S	S	T	Y	-2
	<u>A</u>	<u>W</u>	<u>A</u>	<u>S</u>	<u>Y</u>	<u>R</u>	<u>L</u>	<u>P</u>	<u>T</u>	<u>K</u>	<u>M</u>	<u>M</u>	Q	P	R	H	I	-3	

A reading frame from a start codon to the first stop codon is called an *open* reading frame (underlined above).